PHANTOMWIKI: ON-DEMAND DATASETS FOR REA SONING AND RETRIEVAL EVALUATION

Anonymous authors

004

010 011

012

013

014

015

016

017

018

019

021

Paper under double-blind review

ABSTRACT

High-quality benchmarks are essential for evaluating reasoning and retrieval capabilities of large language models (LLMs). However, curating datasets for this purpose is not a permanent solution as they are prone to data leakage and inflated performance results. To address these challenges, we propose *PhantomWiki*: a pipeline to generate unique, factually consistent document corpora with diverse question-answer pairs. Unlike prior work, *PhantomWiki* is neither a fixed dataset, nor is it based on any existing data. Instead, a new *PhantomWiki* instance is generated on demand for each evaluation. We vary the question difficulty and corpus size to disentangle reasoning and retrieval capabilities, respectively, and find that *PhantomWiki* datasets are surprisingly challenging for frontier LLMs. Thus, we contribute a scalable and data leakage-resistant framework for disentangled evaluation of reasoning, retrieval, and tool-use abilities.

023 024 1 INTRODUCTION

025 Designing agents that can perform complex reasoning while interfacing with a large-scale, dynamic 026 corpus-like Wikipedia-is a long-standing goal in the field of natural language processing (Feldman 027 & El-Yaniv, 2019; Min et al., 2019). Such a goal may be within reach given the impressive capabilities of recent language models, which are all trained on internet-scale data. For example, the ability 029 of LLMs to solve math problems on GSM8K (Cobbe et al., 2021) and mathematical olympiads AlphaProof & AlphaGeometry (2024) could bode well for agents to answer highly quantitative 031 questions. On benchmarks like DROP (Dua et al., 2019) and MMLU (Hendrycks et al., 2020), these LLMs demonstrate advanced reading comprehension and general reasoning capabilities, both 033 necessary for intelligent agents. When augmented with retrievers (Muennighoff et al., 2022) and tools (Patil et al., 2023), LLMs seem to already possess a strong ability for accessing external datastores 034 and knowledge bases. 035

However, it is unclear to what extent these models rely on their internal knowledge, which can easily 037 become outdated, versus their reasoning and retrieval abilities. Consider the example, "What is the 038 date of birth of Wolfgang Amadeus Mozart?". Since this fact is contained within LLMs' pre-training 039 data, asking LLMs this question cannot provide reliable insight on whether the answer was deduced, retrieved or recalled. At the same time, existing approaches that perturb Wikipedia facts (Cohen et al., 040 2024; Meng et al., 2022; Elazar et al., 2021) to construct new question-answer pairs face challenges 041 of ensuring factual consistency across articles. For example, changing Mozart's date of birth to 2025 042 also requires modifying Beethoven's article to erase the fact that Beethoven might have met Mozart 043 in 1787! 044

Reasoning-only benchmarks are not immune to memorization either. On GSM8K, a dataset that contains grade school math problems, Mirzadeh et al. (2024) report that frontier models perform significantly worse with minor or even meaningless alterations to the test data, indicating these models are vulnerable to overfitting at best and exact memorization at worst. To ensure fair comparison, LLMs need to be evaluated in a way that does not depend on any particular dataset instance.

Following this philosophy, we introduce *PhantomWiki*. At the click of a button, *PhantomWiki* generates a fictional universe of characters along with a set of facts. We reflect these facts in a large-scale corpus, mimicking the style of fan-wiki websites. Then we generate question-answer pairs with tunable difficulties, encapsulating the types of multi-hop questions commonly considered in the question-answering (QA) literature.



Figure 1: Evaluating LLM capabilities with *PhantomWiki*. We tune the reasoning and retrieval difficulty by the number of reasoning steps and documents, respectively.

We design *PhantomWiki* to facilitate testing of different aspects of LLMs. In the first setting, the universe is small enough such that all relevant information can fit within the context. Studies such as (Liu et al., 2024) show that LLMs perform poorly in "needle-in-a-haystack" scenarios, where a small but crucial piece of information is embedded within a long document. By adjusting the total context length—determined by *PhantomWiki* universe size—and the quantity of relevant information required for a given question, *PhantomWiki* provides a reliable benchmark for evaluating LLMs' in-context retrieval capabilities.

In the second setting, when dealing with large-scale corpora, LLMs face inherent limitations in processing all available information within their fixed context window. Instead, they must rely on external retrieval methods to access relevant information (Lewis et al., 2020). This setup allows us to decouple and evaluate two key components: the effectiveness of the retriever in identifying and retrieving the most relevant content, and the LLMs' ability to accurately interpret and utilize the retrieved information.

Last but not least, where LLMs are augmented with external tools, effectively integrating reasoning
 and tool-use capabilities becomes essential for solving complex tasks. By adjusting the size of the
 associated text corpus and modulating the reasoning difficulty, *PhantomWiki* serves as a foundation
 for future research on agents that can seamlessly combine reasoning and tool utilization.

091 Our evaluation on PhantomWiki confirms that the proposed tasks present significant challenges for all 092 of the state-of-the-art LLMs that we used. As we show in Figure 1, we observe consistent performance 093 declines across all scenarios as the universe size grows or the number of reasoning steps increases, indicating heightened difficulty in retrieval and reasoning. By breaking down challenges across 094 different dimensions, PhantomWiki enables researchers from various fields to evaluate and refine 095 their methods. Beyond serving as a robust benchmark for LLM performance, *PhantomWiki* provides 096 valuable insights that can guide improvements in retrieval, reasoning, and tool-use capabilities of LLMs for the research community. We will make PhantomWiki code available in a public GitHub 098 repository after the anonymous reviewing period.

100

2 RELATED WORKS

101 102

Agent benchmarks, such as τ -bench (Yao et al., 2024), ToolWoz (Lattimer et al., 2024), Alfworld (Shridhar et al., 2020) and WebArena (Zhou et al., 2023), focus on tasks where the agent is given a binary reward for successful completion of a task (e.g., booking a flight, making a purchase). In this work, we focuses more on tasks where the agent is rewarded for responding to a question with a factually correct answer. (In Section 3, we concretize what we mean by a "fact" and "correctness".) Zhou et al. (2023, Section 3.2) include a category of information-seeking tasks, however these tasks

065 066 067

068 069 070

071

072 073

054

057

060 061 062

125



Figure 2: Overview of the PhantomWiki pipeline.

often require navigation across multiple pages or focus on user-centric content. Yao et al. (2024, Appendix. A) measure task difficulty based on the average success rate of frontier models (e.g., GPT-4). Our work defines a model-agnostic measure of difficulty, which we show provides more meaningful insight into the reasoning and retrieval aspects of LLMs.

135 In the QA domain, existing benchmarks are designed to test whether LLMs are able to reason and 136 use tools. Closer to our work in the space of question-answering agents is the ToolQA benchmark of 137 Zhuang et al. (2023). They introduce a framework to construct question-answer pairs from databases and documents by first generating question templates using LLMs, then filtering for high-quality 138 templates, and finally deriving ground-truth answers by writing corresponding Python programs 139 for each question template. Zhuang et al. (2023, Tab. 1) construct two pure-text datasets: SciREX 140 with 438 documents and 5 question templates, and Agenda with 10k event entries and 10 question 141 templates. In constrast, PhantomWiki generates instances at much larger scale with 50 question 142 templates and 1 million documents. 143

For generating factually-grounded answers, retrieval augmented generation (RAG) has emerged as 144 the predominant paradigm (Lewis et al., 2020; Karpukhin et al., 2020; Guu et al., 2020). However, 145 evaluating RAG systems is notoriously difficult, leading to a flourishing of retrieval benchmarks 146 (Petroni et al., 2020; Saad-Falcon et al., 2023; Jin et al., 2024; Hsia et al., 2024; Mao et al., 2024; Rau 147 et al., 2024). A key pain-point of RAG is handling questions that involve multi-hop reasoning. This 148 motivated Tang & Yang (2024) to design the MultiHop-RAG dataset with synthetically generated 149 questions and Su et al. (2024) to curate a dataset of question-answer pairs that requires intensive rea-150 soning to retrieve relevant documents. Importantly, none of these benchmarks creates the underlying 151 corpus, a limitation which we bridge in this work.

152 Logical reasoning tasks have become central to LLM evaluation and have garnered significant 153 attention in recent time (Zhu et al., 2023). Many existing benchmarks do not disentangle the evaluation 154 on logical reasoning with other abilities such as natural language inference and commonsense 155 reasoning (Sakaguchi et al., 2021; Zellers et al., 2019; Sprague et al., 2023). A line of works focuses 156 on the synthesis of datasets containing a variety of logic reasoning tasks (Tafjord et al., 2020; Saparov 157 & He, 2022; Liu et al., 2020; Han et al., 2022; Weston et al., 2015). Closer to our work, Sinha et al. 158 (2019) construct short stories about individuals related through a family graph and ask questions 159 about their kinship relationships to benchmark the inductive reasoning capabilities. However, theirs is distinct from our work in that all relevant information for a specific question is centralized in a single 160 article; PhantomWiki requires that the relevant information first be retrieved from a large-scale 161

corpus.

162 PhantomWiki CONSTRUCTION 3

163 164

PhantomWiki is at its core an on-demand random generator of fictional worlds. Similarly to the wiki 165 hosting services popular in film, video games and literature¹, we represent these fictional worlds 166 through Wikipedia-like biographical entries about their characters. We then test the model's retrieval 167 skills and its understanding of the fictional world through an accompanying set of automatically 168 generated question-answer pairs.

169 170

3.1 GENERATING A PhantomWiki UNIVERSE

171 The first part of the *PhantomWiki* pipeline generates a random universe of n characters as well as the 172 document corpus describing it, as illustrated in Figure 2, (1-2). 173

174 Generating Characters. Each character in a PhantomWiki universe is described through its social 175 relationships and personal facts (Figure 2, (1)). For the social relationships, we first generate family 176 trees, following the family tree generator from Hohenecker & Lukasiewicz (2020). We iteratively pick a person and generate their parent or child based on various constraints², until the user-specified 177 universe size of n people is reached. The user can also specify other hyperparameters like the number 178 of trees, their maximal depth, and the maximal number of offspring for each person. In addition to the 179 family trees, we generate a friendship graph using the Erdős–Rényi model (making two people friends 180 with some fixed probability, typically controlled by the desired average number of friendships.) 181

Generating Facts. Next, we generate personal facts for each person in the *PhantomWiki* universe. 182 Names are assigned during the family generation procedure, with the first name sampled based on 183 the character's gender and the surname based on the family tree, resulting in 15M full names in total³. We also add dates of birth in a way that is consistent with the existing family relations, and assign 185 each person a job and a hobby that we uniformly sample from over 300 and 600 options respectively.

187 Generating Articles. Given all relevant facts for each person, we convert them into articles using pre-defined templates, e.g. "The job of David is a farmer. The hobby of David is birdwatching." 188 (see Figure 2, (2)). This construction conveys the necessary information while keeping the articles 189 short (about 160 tokens on average). While it is possible to extend the article generation process to 190 LLM-based methods (see e.g. Shao et al. 2024), this poses the challenge of guaranteeing factual 191 correctness without additional costs and external supervision. This has been supported by our 192 preliminary experiments on article generation using Llama-3.3-70B, where we observed factual errors 193 in the resulting articles; therefore we do not use LLMs and rely entirely on templates. The articles 194 are the only component of *PhantomWiki* available to the model during its evaluation.

195 196 197

207 208 209

213

3.2 GENERATING QUESTION-ANSWER PAIRS

In the second half of the *PhantomWiki* pipeline, we generate a set of questions with verifiable 199 answers, as shown in Figure 2, (3-4).

200 Generating Questions. We implement automatic question generation through a context-free grammar 201 (CFG, Hopcroft et al. 2001) of question templates, which we then use to sample complete questions. 202 For example, the question template "Who is the <relation> of <name>?" can be used to sample the 203 question "Who is the friend of David?" (see Figure 2, (3)). The main advantage of using a CFG is 204 that it efficiently and systematically obtains all possible compositions of questions for some recursion 205 limit d. For instance, the following subset of our context-free grammar: 206

$$R \rightarrow \text{the } < relation > \text{of } R'$$

$$R' \rightarrow R \mid < name >$$

 $S \rightarrow$ Who is R?

210 can lead to questions ranging from "Who is the friend of David?" to "Who is the nephew of the friend 211 of the brother of David?" as d increases. In addition to these nested compositions, our CFG also 212

¹For example, see stardewvalley.fandom.com or harrypotter.fandom.com.

²For example, the number of offspring of a person has to be smaller than some threshold, parents of the 214 people at the maximal tree level will not be generated, etc. 215

³We use unique names in our experiments, but *PhantomWiki* also supports repeated names.

supports questions about personal attributes (e.g. "Who is the person whose hobby is birdwatching?"),
aggregation questions ("How many brothers does David have?"), and combinations of all three ("How many friends does the brother of the person whose hobby is birdwatching have?") (For the full CFG
see Appendix B.)

220 **Generating Answers.** To ensure that the answers to the sampled questions are verifiably correct, we 221 represent our generated universe in Prolog, a logic programming language (Sterling & Shapiro, 1994). 222 Each Prolog program consists of a set of facts known about the world such as hobby ("David", 223 "birdwatching"), and a set of rules defining how facts are related to each other, such as 224 nephew (X, Y) :- sibling (X, A), son (A, Y). The Prolog program uses these facts 225 and rules to deduce the exhaustive set of answers to its queries (i.e., the CFG-generated questions). 226 For example, a question "Who is the nephew of the friend of the person whose hobby is birdwatching?" corresponds to the three-statement Prolog query ?- nephew(X2, Y), friend(X1, X2), 227 hobby (X1, "birdwatching"), which returns all people satisfying these constraints in the 228 PhantomWiki universe (see Figure 2 (4)). 229

To construct the Prolog queries automatically, we modify the CFG algorithm to generate both the
 question and query templates in parallel. We note, however, that the queries are separate from the final
 PhantomWiki corpus and question-answer pairs, and the answers returned by the Prolog program
 should be held out as part of the evaluation procedure.

234 235 236

3.3 PhantomWiki COMPLEXITY

The goal of *PhantomWiki* is to generate memorization-resistant evaluation datasets that are challenging in both reasoning and retrieval aspects. In this section, we discuss our conceptual and practical design choices that help us achieve this goal.

240 **Universe Space Complexity.** To ensure that our evaluation with *PhantomWiki* is memorization 241 and data leakage-resistant, we first show that the space of possible universes is sufficiently large to generate enough unique instances. Observe that the number of possible friendship assignments 242 grows at the rate of $\Theta(2^{n^2})$ (Flajolet & Sedgewick, 2009, Ex. II.5) as the number of individuals n in 243 the universe increases. Similarly, assuming each individual is assigned one fact from each category 244 (job, hobby, etc.), the number of possible fact assignments grows at the rate $\Theta(c^n)$, where c is the 245 total number of choices across the categories. *PhantomWiki* thus samples a corpus from $\Theta(2^{n^2}c^n)$ 246 possible universes, which leads to diverse datasets optimal for data leakage-resistant evaluation. We 247 note that as future work *PhantomWiki* could be extended to increase this diversity, e.g. by adding a 248 temporal dimension of events. 249

250 **Reasoning Complexity.** The CFG enables us to recursively compose templates that lead to complex 251 reasoning questions. Observe that our CFG in Appendix **B** produces $\Theta(d)$ question templates as the 252 recursion limit d increases. Moreover, we can increase the difficulty of each template by increasing the number of *reasoning steps*. For example, substituting *<relation>* with *nephew* in a template adds two 253 reasoning steps (nephew (X, Y) :- sibling (X, A), son (A, Y)), since PhantomWiki 254 articles only contain immediate family relationships like *sibling* and *son*. In contrast, substituting 255 <relation> with second cousin would lead to five reasoning steps. As we will show in Section 4, 256 *PhantomWiki* questions are sufficiently complex to evaluate reasoning capabilities of state-of-the-art 257 LLMs. We further note that *PhantomWiki*'s CFG can be easily extended to support more question 258 types like comparison and multiple-constraint questions. 259

Retrieval Complexity. To assess a model's retrieval capabilities, we increase the universe size n260 so that the document corpus exceeds the model's context length—this makes a retriever necessary 261 to answer questions correctly. For state-of-the-art LLMs with a context length of 128K, such as 262 OpenAI's GPT-40 and Meta's Llama-3.3-70B, this corresponds to *PhantomWiki* universes of $n \gtrsim 1$ K. 263 This increases to $n \gtrsim 3K$ for Google's Gemini-1.5-Flash with context length 1M. Further scaling n264 leads to further increase in retrieval difficulty. In Table 1, we show that *PhantomWiki* is well-suited 265 for generating universes of this size on standard CPU hardware: generating 10-hop questions for 266 size n = 100K—well beyond any existing LLM's context length—takes just 6 minutes on 8 Intel 267 Cascade Lake CPU cores. Moreover, we can conveniently generate instances of n = 1M, which is on 268 the scale of Wikipedia's corpus of 2 million biographical entries⁴.

⁴WikiProject Biography, as of January 30, 2025.

Table 1: Runtimes of generating a *PhantomWiki* instance for different universe sizes n.

n	Runtime
10^{2}	1 second
10^{3}	3 seconds
10^{4}	21 seconds
10^{5}	6 minutes
10^{6}	4 hours

Table 2: F1 scores for various LLMs and prompting techniques. We report mean \pm standard error across 3 dataset generation seeds (except for GPT-40 due to cost constraints), and indicate the highest F1 score for each n in bold. In-Context prompting is infeasible for n = 5K as the corpus cannot be fully included in the context.

284		Mada1	In-Context		RAG		Agentic
285	n	Widdel	ZEROSHOT	СоТ	ZEROSHOT-	CoT-RAG	REACT
286					RAG		
207	50	DeepSeek-R1-32B	42.42 ± 1.69	52.42 ± 2.64	16.87 ± 0.98	16.10 ± 1.10	_
289	50	GPT-40	27.20	50.66	20.54	14.14	38.70
200	50	Gemini-1.5-Flash	28.49 ± 1.15	34.61 ± 2.41	19.88 ± 2.05	13.35 ± 0.66	30.92 ± 1.41
290	50	Llama-3.3-70B	25.64 ± 0.56	48.05 ± 1.95	17.55 ± 2.20	20.01 ± 1.81	35.83 ± 1.00
292	500	DeepSeek-R1-32B	18.33 ± 2.33	19.65 ± 3.00	12.08 ± 1.07	9.64 ± 0.46	
293	500	GPT-40	16.76	41.02	13.56	7.25	37.39
294	500	Gemini-1.5-Flash	17.39 ± 1.45	25.17 ± 1.77	11.66 ± 0.34	7.94 ± 0.19	26.99 ± 1.84
295	500	Llama-3.3-70B	11.59 ± 1.19	28.07 ± 0.37	10.89 ± 0.58	11.54 ± 1.05	35.56 ± 0.49
296	5000	DeepSeek-R1-32B			8.29 ± 0.18	7.96 ± 0.38	
297	5000	GPT-40	Max o	context	10.12	6.96	36.85
298	5000	Gemini-1.5-Flash	exce	eded	8.60 ± 0.31	5.26 ± 0.36	23.47 ± 1.53
299	5000	Llama-3.3-70B			7.57 ± 0.52	8.67 ± 0.18	30.89 ± 2.24

EXPERIMENTAL VALIDATION

We evaluate reasoning and retrieval capabilities of several frontier LLMs using PhantomWiki, by decomposing their performance over questions of varying difficulty and universes of varying sizes.

4.1 EVALUATION SETUP

We generate *PhantomWiki* instances with n ranging from 50 to 10K—a universe size for which the total length of articles exceed the LLM context length. For the evaluation, only the articles (not the Prolog database or the generated graphs) will be provided to the LLMs. To ensure that our findings are not tied to any specific PhantomWiki instance, we use 3 random dataset seeds for each configuration. Creating PhantomWiki instances with different random seeds leads to entirely different combinations of names, relations, and personal facts. In each instance, we generate question templates with maximum recursion depth d = 20, for a total of 50 templates. We sample 10 questions for each template, yielding a total of 500 questions per PhantomWiki instance. As shown in Figures 5 and 6 (Appendix B), these questions have varying difficulty and number of answers. Accordingly, we prompt the LLMs to predict all answers as a comma-separated list and measure correctness with the answer-level F1 score.

4.2 MODELS AND PROMPTING TECHNIOUES

We test both open- and closed-source LLMs, namely OpenAI's GPT-40 (Hurst et al., 2024), Google's Gemini-1.5-Flash (Gemini Team, Google, 2024), and the instruction-tuned version of Meta's Llama-3.3-70B model (Dubey et al., 2024). We also evaluate DeepSeekAI's DeepSeek-R1-32B (Guo et al.,

2025) (distilled with Qwen-2.5-32B (Yang et al., 2024)), which is an open-weights LLM trained on
 reasoning trace datasets. We prompt each LLM with the following techniques, broadly grouped in
 three ways:

 In-Context Prompting. This technique includes the whole document corpus as part of the prompt. We use this type of prompting in conjunction with two strategies: ZEROSHOT—where the document corpus is immediately followed by the question—and Chain-of-Thought (COT) prompting (Wei et al., 2022), where we additionally include some examples on how step-by-step reasoning could lead to the correct answer. We include these prompts in Appendix C.2, as well as modifications to the ZEROSHOT strategy for DeepSeek-R1-32B.

RAG Prompting. This setting augments generation with a pre-trained neural retriever (Lewis et al., 2020). We implement this by first searching for the 4 most relevant documents to the posed question based on UAE-LARGE-V1 embeddings. Next, we incorporate these retrieved documents into the model's prompt. Finally, we add in the same ZEROSHOT and COT prompts as In-Context Prompting. We document details on our retrieval algorithm in Appendix C.3.

Agentic Prompting. REACT (Yao et al., 2022) is a prompting technique that enables LLMs to interleave reasoning steps with tool interactions, to solve complex tasks. For *PhantomWiki* QA task, the LLMs are provided with keyword-based tools RetrieveArticle and Search to retrieve relevant documents (see Appendix C.6 for tool details). These settings materialize the limitations of in-context prompting and necessitate the use of advanced RAG prompting and agentic prompting approaches.

In the CoT and REACT prompts, we include 10 QA exemplars and hand-written reasoning traces. We choose these exemplars from a dataset instance of size 25 that is not used for evaluation. In REACT, we limit LLMs to interact with the text corpus for up to 50 steps, which is sufficient to answer almost all questions in *PhantomWiki* instances.

We cap all LLM outputs to 4096 tokens and use greedy decoding (temperature = 0). For DeepSeek-R1-32B, we use temperature = 0.6 and top-p = 0.95 in accordance with the evaluation setup in Guo et al. (2025, Sec. 3). We refer the reader to Appendix C for full prompt templates and implementation details.

4.3 DISCUSSION

In Table 2, we report the mean F1 score across various universe sizes, LLMs, and prompting techniques. We first average F1 scores over all questions in a *PhantomWiki* instance, then compute the mean and standard error across the dataset generation seeds.

358 We first consider the small-universe setting (n = 50) in Table 2, which corresponds to roughly 16K 359 tokens for the LLMs we test. In-Context prompting techniques outperform other techniques: COT 360 with GPT-40 attains the highest performance, followed by ZEROSHOT with DeepSeek-R1-32B. Next, 361 we consider the setting of medium universes (n = 500). Here the full document corpus can still 362 be included in all LLMs' contexts, but we find that ZEROSHOT performs poorly for all LLMs, and 363 DeepSeek-R1-32B especially struggles. F1 scores of CoT for all LLMs degrade as well compared to 364 n = 50, but not worse than REACT. Finally, in the setting of large universes (n = 5000), none of 365 the LLMs we evaluate can accommodate the full document corpus. in-context prompting techniques 366 are no longer viable, and we must rely on RAG prompting and agentic prompting. RAG prompting attain poor F1 scores because the retriever fails to retrieve documents relevant for answering complex 367 questions. On the other hand, agentic prompting technique shines in comparison to other techniques, 368 indicating that LLMs are better suited to dynamically retrieve documents while reasoning on a 369 question. 370

371

373

353

354

372 5 EVALUATING REASONING

To isolate LLM reasoning capabilities with *PhantomWiki*, we investigate model performance on small universes (n = 50) in Figure 3. Note that contexts of all LLMs can fully include small universe document corpora. Each *PhantomWiki* dataset contains questions covering a wide range of difficulty. We evaluate three approaches: in-context prompting, RAG prompting, and agentic prompting. For each we plot the F1 scores as a function of question difficulty, as measured by the number of *reasoning* 391

392

393

394

395 396 397



Figure 3: F1 scores as a function of question difficulty, measured by *reasoning steps*. We plot LLM performance on universe size n = 50, and report F1 scores averaged over 3 generation seeds. Increasing question difficulty in *PhantomWiki* reveals a clear decline across all state-of-the-art LLMs and prompting techniques, showing their struggle with reasoning.

steps necessary to answer the question. As mentioned in Section 3.3, this is determined by the type of
 question templates and the sampled relationships. For all LLMs and prompting techniques, we verify
 empirically that questions with larger reasoning steps are indeed more challenging to answer.
 By allowing question difficulty to be adjusted, *PhantomWiki* serves as a foundational benchmark for
 evaluating reasoning capabilities in language models.

ZEROSHOT performance declines sharply as the number of reasoning steps increases for all LLMs,
 except for DeepSeek-R1-32B, which deteriorates more gradually. LLMs perform better with COT
 than with ZEROSHOT, but each additional reasoning step remains increasingly challenging. This
 suggests that even in the absence of retrieval constraints, LLMs struggle to navigate logical reasoning
 sequences.

RAG prompting techniques (ZEROSHOT-RAG and COT-RAG) stunt reasoning performance across
 the board—F1 scores are near zero on questions with 5 or more reasoning steps as opposed to 15
 steps for in-context prompting. We attribute this to a core problem with RAG prompting: retrieving
 documents in the initial prompt before starting to answer the question, as opposed to reasoning
 through the question and retrieving documents dynamically.

413 We find that RAG prompting techniques can only answer questions that require a single reasoning step, 414 like Who is the friend of David?. On the other hand, answering questions that require information from 415 multiple reasoning steps is extremely challenging for ZEROSHOT-RAG and COT-RAG. To illustrate, 416 consider the question Who is the nephew of the friend of David? Answering this question requires 417 retrieving David's document first and then retrieving their friend's document to find the nephew. Since RAG prompting techniques retrieve documents *only once* by matching vector embeddings of 418 questions and documents, they are unlikely to retrieve all necessary documents required to answer 419 such questions. 420

Finally, the agentic prompting technique REACT allows LLMs to avoid the steep performance drop as seen in RAG prompting. On given a question, REACT prompting requires LLMs to retrieve documents dynamically in a conversation and justify why they are relevant. Concretely, before using a tool (RetrieveArticle or Search) in a conversation turn, the LLM is asked to describe how the tool will help using a "Thought" step (Yao et al., 2022), analogous to the COT prompting approach. This approach shows promise in answering questions correctly. Even so, REACT struggles as the question difficulty increases.

Figure 3 thus decomposes LLM performance along the lines of reasoning capabilities. It reveals
that all in-context prompting and agentic prompting achieve near-perfect F1 scores on low-difficulty
questions. Therefore, the stratification between them in Table 2 can be attributed to varying performance on high difficulty questions. To further isolate the impact of question difficulty, in Figure 7 we plot F1 scores as a function of reasoning steps for questions with only one solution.



Figure 4: **F1 scores as a function of universe size** n. We evaluate LLM performance on questions with ≤ 10 reasoning steps, and report F1 scores averaged over 3 dataset generation seeds. As we increase universe size in *PhantomWiki*, F1 scores for all LLMs and prompting techniques deteriorate, highlighting that they struggle at retrieving relevant documents.

6 EVALUATING RETRIEVAL

Next, to evaluate LLM retrieval capabilities, we use *PhantomWiki* to contrast two settings: (1) small universes where the document corpus can comfortably fit in LLM context, and (2) large universes where the full corpus exceeds context lengths. To this end, we increase the universe size up to n = 10K, which corresponds to document corpora well beyond the context lengths of state-of-the-art LLMs, and display the results in Figure 4.

For very small universes, COT usually outperforms ZEROSHOT for all LLMs except DeepSeek-R1-32B. However, F1 scores noticeably worsen as more documents are included in models' contexts, with DeepSeek-R1-32B suffering a dramatic performance drop. This analysis regime indicates that state-of-the-art LLMs struggle at in-context retrieval for complex question-answering tasks.

At the large universe scale, in-context prompting techniques become nonviable as the document corpus exceeds model context lengths. Therefore the use of out-of-context retrieval, such as RAG prompting and agentic prompting techniques, is necessary for obtaining the answers. Here we observe that RAG prompting techniques, which select relevant documents for the question using vector embeddings, deliver poor F1 scores for all universe sizes—the performance only deteriorates with increasing universe size. Agentic prompting techniques like REACT show immense promise by avoiding a steep downward trend. This suggests that **agentic workflows can be effective in dynamically retrieving documents at scale**.

471 472

445

446

447

448

449 450 451

452 453 454

455

456

457

458

7 CONCLUSION AND FUTURE WORK

- 473
- 475

We introduce *PhantomWiki*: a benchmarking framework to evaluate reasoning and retrieval capabilities of language models. As we increase the question complexity and universe size, we observe that current state-of-the-art LLMs struggle in both reasoning and retrieval aspects. *PhantomWiki* is scalable and memorization-resistant, hence well-suited to evaluate future generations of language models.

Our work brings forth several research directions. Noting how we generate document corpora and questions, *PhantomWiki* is resistant to data contamination. We leave to future work to empirically test this claim, and develop theory to formally prove that our benchmark is memorization-resistant. In this work we focus on question-answering over text corpora. We hope to extend *PhantomWiki* for other knowledge bases and modalities such as vision and audio, enabling analogous evaluation suites for multimodal models.

486	REFERENCES
487	

494

501

511

523

524

525

488	AlphaProof & AlphaGeometry. AI achieves silver-medal standard solving International Mathematical
489	Olympiad problems. https://deepmind.google/discover/blog/ai-solves-imo-problems-at-silver-
490	medal-level/, 2024. Accessed: 2025-01-25.

- 491 Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Mark Chen, Heewoo Jun, Lukasz Kaiser, 492 Matthias Plappert, Jerry Tworek, Jacob Hilton, Reiichiro Nakano, et al. Training verifiers to solve 493 math word problems. arXiv preprint arXiv:2110.14168, 2021.
- Roi Cohen, Eden Biran, Ori Yoran, Amir Globerson, and Mor Geva. Evaluating the ripple effects 495 of knowledge editing in language models. Transactions of the Association for Computational 496 Linguistics, 12:283-298, 2024. 497
- 498 Dheeru Dua, Yizhong Wang, Pradeep Dasigi, Gabriel Stanovsky, Sameer Singh, and Matt Gardner. 499 DROP: A reading comprehension benchmark requiring discrete reasoning over paragraphs. arXiv 500 preprint arXiv:1903.00161, 2019.
- Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha 502 Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, et al. The Llama 3 herd of models. arXiv preprint arXiv:2407.21783, 2024. 504
- 505 Yanai Elazar, Nora Kassner, Shauli Ravfogel, Abhilasha Ravichander, Eduard Hovy, Hinrich Schütze, 506 and Yoav Goldberg. Measuring and improving consistency in pretrained language models. Trans-507 actions of the Association for Computational Linguistics, 9:1012–1031, 2021.
- 508 Yair Feldman and Ran El-Yaniv. Multi-hop paragraph retrieval for open-domain question answering. 509 arXiv preprint arXiv:1906.06606, 2019. 510
- P. Flajolet and R. Sedgewick. Analytic Combinatorics. Cambridge University Press, 2009. ISBN 512 9781139477161. URL https://books.google.com/books?id=0h-4QcA1c1QC. 513
- Gemini Team, Google. Gemini 1.5: Unlocking multimodal understanding across millions of tokens 514 of context. arXiv preprint arXiv:2403.05530, 2024. 515
- 516 Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song, Ruoyu Zhang, Runxin Xu, Qihao Zhu, 517 Shirong Ma, Peiyi Wang, Xiao Bi, et al. DeepSeek-R1: Incentivizing reasoning capability in 518 LLMs via reinforcement learning. arXiv preprint arXiv:2501.12948, 2025. 519
- Kelvin Guu, Kenton Lee, Zora Tung, Panupong Pasupat, and Mingwei Chang. Retrieval augmented 520 language model pre-training. In International conference on machine learning, pp. 3929–3938. 521 PMLR, 2020. 522
 - Simeng Han, Hailey Schoelkopf, Yilun Zhao, Zhenting Qi, Martin Riddell, Wenfei Zhou, James Coady, David Peng, Yujie Qiao, Luke Benson, et al. Folio: Natural language reasoning with first-order logic. arXiv preprint arXiv:2209.00840, 2022.
- 526 Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou, Mantas Mazeika, Dawn Song, and 527 Jacob Steinhardt. Measuring massive multitask language understanding. arXiv preprint 528 arXiv:2009.03300, 2020. 529
- 530 Patrick Hohenecker and Thomas Lukasiewicz. Ontology reasoning with deep neural networks. 531 Journal of Artificial Intelligence Research, 68:503–540, 2020.
- John E Hopcroft, Rajeev Motwani, and Jeffrey D Ullman. Introduction to automata theory, languages, 533 and computation. Acm Sigact News, 32(1):60-65, 2001. 534
- 535 Jennifer Hsia, Afreen Shaikh, Zhiruo Wang, and Graham Neubig. RAGGED: Towards informed 536 design of retrieval augmented generation systems. arXiv preprint arXiv:2403.09040, 2024. 537
- Aaron Hurst, Adam Lerer, Adam P Goucher, Adam Perelman, Aditya Ramesh, Aidan Clark, AJ Os-538 trow, Akila Welihinda, Alan Hayes, Alec Radford, et al. Gpt-4o system card. arXiv preprint arXiv:2410.21276, 2024.

540	Jiaije Jin, Yutao Zhu, Xinyu Yang, Chenghao Zhang, and Zhicheng Dou. FlashRAG: A modular
541	toolkit for efficient retrieval-augmented generation research arXiv prantiat arXiv:2405.13576
E40	tookit for enclent retreval-augmented generation research. <i>urxiv preprint urxiv.2405.15570</i> ,
342	2024.
543	

- Vladimir Karpukhin, Barlas Oğuz, Sewon Min, Patrick Lewis, Ledell Wu, Sergey Edunov, Danqi
 Chen, and Wen-tau Yih. Dense passage retrieval for open-domain question answering. *arXiv preprint arXiv:2004.04906*, 2020.
- Barrett Martin Lattimer, Varun Gangal, Ryan McDonald, and Yi Yang. Sparse rewards can self-train dialogue agents. *arXiv preprint arXiv:2409.04617*, 2024.
- Patrick Lewis, Ethan Perez, Aleksandra Piktus, Fabio Petroni, Vladimir Karpukhin, Naman Goyal, Heinrich Küttler, Mike Lewis, Wen-tau Yih, Tim Rocktäschel, et al. Retrieval-augmented generation for knowledge-intensive nlp tasks. *Advances in Neural Information Processing Systems*, 33: 9459–9474, 2020.
- Jian Liu, Leyang Cui, Hanmeng Liu, Dandan Huang, Yile Wang, and Yue Zhang. Logiqa: A
 challenge dataset for machine reading comprehension with logical reasoning. *arXiv preprint arXiv:2007.08124*, 2020.
- Nelson F Liu, Kevin Lin, John Hewitt, Ashwin Paranjape, Michele Bevilacqua, Fabio Petroni, and
 Percy Liang. Lost in the middle: How language models use long contexts. *Transactions of the Association for Computational Linguistics*, 12:157–173, 2024.
- Qianren Mao, Yangyifei Luo, Jinlong Zhang, Hanwen Hao, Zhilong Cao, Xiaolong Wang, Xiao
 Guan, Zhenting Huang, Weifeng Jiang, Shuyu Guo, et al. XRAG: eXamining the core benchmarking foundational components in advanced retrieval-augmented generation. *arXiv preprint arXiv:2412.15529*, 2024.
- Kevin Meng, David Bau, Alex Andonian, and Yonatan Belinkov. Locating and editing factual associations in gpt. *Advances in Neural Information Processing Systems*, 35:17359–17372, 2022.
- Sewon Min, Danqi Chen, Hannaneh Hajishirzi, and Luke Zettlemoyer. A discrete hard em approach
 for weakly supervised question answering. *arXiv preprint arXiv:1909.04849*, 2019.
- Iman Mirzadeh, Keivan Alizadeh, Hooman Shahrokhi, Oncel Tuzel, Samy Bengio, and Mehrdad
 Farajtabar. GSM-Symbolic: Understanding the limitations of mathematical reasoning in large
 language models. *arXiv preprint arXiv:2410.05229*, 2024.
- Niklas Muennighoff, Nouamane Tazi, Loïc Magne, and Nils Reimers. MTEB: Massive text embed ding benchmark. *arXiv preprint arXiv:2210.07316*, 2022.
- Shishir G Patil, Tianjun Zhang, Xin Wang, and Joseph E Gonzalez. Gorilla: Large language model connected with massive APIs. *arXiv preprint arXiv:2305.15334*, 2023.
- Fabio Petroni, Aleksandra Piktus, Angela Fan, Patrick Lewis, Majid Yazdani, Nicola De Cao, James
 Thorne, Yacine Jernite, Vladimir Karpukhin, Jean Maillard, et al. Kilt: a benchmark for knowledge
 intensive language tasks. *arXiv preprint arXiv:2009.02252*, 2020.
- David Rau, Hervé Déjean, Nadezhda Chirkova, Thibault Formal, Shuai Wang, Vassilina Nikoulina, and Stéphane Clinchant. BERGEN: A benchmarking library for retrieval-augmented generation. *arXiv preprint arXiv:2407.01102*, 2024.
- Jon Saad-Falcon, Omar Khattab, Christopher Potts, and Matei Zaharia. ARES: An automated evaluation framework for retrieval-augmented generation systems. *arXiv preprint arXiv:2311.09476*, 2023.
- Keisuke Sakaguchi, Ronan Le Bras, Chandra Bhagavatula, and Yejin Choi. WinoGrande: An adversarial Winograd schema challenge at scale. *Communications of the ACM*, 64(9):99–106, 2021.
- 593 Abulhair Saparov and He He. Language models are greedy reasoners: A systematic formal analysis of chain-of-thought. *arXiv preprint arXiv:2210.01240*, 2022.

594 Yijia Shao, Yucheng Jiang, Theodore A Kanell, Peter Xu, Omar Khattab, and Monica S Lam. 595 Assisting in writing Wikipedia-like articles from scratch with large language models. arXiv 596 preprint arXiv:2402.14207, 2024. 597 Noah Shinn, Federico Cassano, Ashwin Gopinath, Karthik Narasimhan, and Shunyu Yao. Reflexion: 598 Language agents with verbal reinforcement learning. Advances in Neural Information Processing Systems, 36, 2024. 600 601 Mohit Shridhar, Xingdi Yuan, Marc-Alexandre Côté, Yonatan Bisk, Adam Trischler, and Matthew 602 Hausknecht. Alfworld: Aligning text and embodied environments for interactive learning. arXiv 603 preprint arXiv:2010.03768, 2020. 604 Koustuv Sinha, Shagun Sodhani, Jin Dong, Joelle Pineau, and William L Hamilton. CLUTRR: A 605 diagnostic benchmark for inductive reasoning from text. arXiv preprint arXiv:1908.06177, 2019. 606 607 Zayne Sprague, Xi Ye, Kaj Bostrom, Swarat Chaudhuri, and Greg Durrett. MuSR: Testing the limits 608 of chain-of-thought with multistep soft reasoning. arXiv preprint arXiv:2310.16049, 2023. 609 Leon Sterling and Ehud Y Shapiro. The art of Prolog: advanced programming techniques. MIT 610 press, 1994. 611 612 Hongjin Su, Howard Yen, Mengzhou Xia, Weijia Shi, Niklas Muennighoff, Han-yu Wang, Haisu Liu, 613 Quan Shi, Zachary S Siegel, Michael Tang, et al. BRIGHT: A realistic and challenging benchmark for reasoning-intensive retrieval. arXiv preprint arXiv:2407.12883, 2024. 614 615 Oyvind Tafjord, Bhavana Dalvi Mishra, and Peter Clark. Proofwriter: Generating implications, 616 proofs, and abductive statements over natural language. arXiv preprint arXiv:2012.13048, 2020. 617 Yixuan Tang and Yi Yang. MultiHop-RAG: Benchmarking retrieval-augmented generation for 618 multi-hop queries. arXiv preprint arXiv:2401.15391, 2024. 619 620 Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny 621 Zhou, et al. Chain-of-thought prompting elicits reasoning in large language models. Advances in 622 neural information processing systems, 35:24824–24837, 2022. 623 Jason Weston, Antoine Bordes, Sumit Chopra, Alexander M Rush, Bart Van Merriënboer, Armand 624 Joulin, and Tomas Mikolov. Towards AI-complete question answering: A set of prerequisite toy 625 tasks. arXiv preprint arXiv:1502.05698, 2015. 626 627 An Yang, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chengyuan Li, 628 Dayiheng Liu, Fei Huang, Haoran Wei, et al. Qwen 2.5 Technical Report. arXiv preprint 629 arXiv:2412.15115, 2024. 630 Shunyu Yao, Jeffrey Zhao, Dian Yu, Nan Du, Izhak Shafran, Karthik Narasimhan, and Yuan Cao. 631 ReAct: Synergizing reasoning and acting in language models. arXiv preprint arXiv:2210.03629, 632 2022. 633 634 Shunyu Yao, Noah Shinn, Pedram Razavi, and Karthik Narasimhan. τ -bench: A benchmark for 635 tool-agent-user interaction in real-world domains. arXiv preprint arXiv:2406.12045, 2024. 636 Rowan Zellers, Ari Holtzman, Yonatan Bisk, Ali Farhadi, and Yejin Choi. HellaSwag: Can a machine 637 really finish your sentence? arXiv preprint arXiv:1905.07830, 2019. 638 639 Shuyan Zhou, Frank F Xu, Hao Zhu, Xuhui Zhou, Robert Lo, Abishek Sridhar, Xianyi Cheng, 640 Tianyue Ou, Yonatan Bisk, Daniel Fried, et al. WebArena: A realistic web environment for building autonomous agents. arXiv preprint arXiv:2307.13854, 2023. 641 642 Zhaocheng Zhu, Yuan Xue, Xinyun Chen, Denny Zhou, Jian Tang, Dale Schuurmans, and Hanjun 643 Dai. Large language models can learn rules. arXiv preprint arXiv:2310.07064, 2023. 644 645 Yuchen Zhuang, Yue Yu, Kuan Wang, Haotian Sun, and Chao Zhang. ToolQA: A dataset for LLM 646 question answering with external tools. Advances in Neural Information Processing Systems, 36: 50117-50143, 2023. 647

648 BACKGROUND А 649

653

654

655

657 658

659

660 661

662 663

664 665

666

676

677 678

679

691

692

650 A.1 CONTEXT-FREE GRAMMARS 651

652 Context-free grammar (CFG) is a type of formal grammar where the productions rules govern how to generate text from non-terminals and terminals. A context-free grammar is defined by $G = (V, \Sigma, R, S)$ where V and Σ denotes nonterminal and terminal respectively. R is a finite relation in $V \times (V \cup \Sigma)^*$ which specifies the production rules of the grammar. $S \in V$ is the start symbol. A production rule in R has the form 656

$$\alpha \to \beta \tag{1}$$

where $\alpha \in V, \beta \in (V \cup \Sigma)^*$. It is conventional to list all rules with the same left-hand side on the same line and separate the right-hand side with "|" like $\alpha \rightarrow \beta_1 | \beta_2$.

QUESTION TEMPLATE GENERATION В

B.1 CONTEXT-FREE GRAMMAR

We use the following CFG to generate question templates:

```
667
          S -> Who is R? | What is A ? | How many RN_p does R_c have ?
          R -> the RN of R_c | the person whose AN is AV
668
          R_c -> R | N
669
          A -> the AN of R
670
          RN -> <relation>
671
          RN_p -> <relation_plural>
672
          AN -> <attribute_name>
673
          AV -> <attribute_value>
674
          N -> <name>
675
```

B.2 CFG-GENERATED QUESTION TEMPLATES

Our CFG produces the following 50 question templates at recursion limit d = 20. Note how the recursive production rule $R_c \rightarrow R \mid N$ leads to chained productions.

```
680
      1. Who is the <relation>_3 of the <relation>_5 of the <relation>_7
681
          of the <relation>_9 of the <relation>_11 of the <relation>_13
682
          of the <relation>_15 of the <relation>_17 of the person whose
683
          <attribute_name>_19 is <attribute_value>_19?
684
      2. Who is the <relation>_3 of the <relation>_5 of the <relation>_7
685
          of the <relation>_9 of the <relation>_11 of the <relation>_13
686
```

```
of the <relation>_15 of the <relation>_17 of <name>_18?
      3. Who is the <relation> 3 of the <relation> 5 of the <relation> 7
688
          of the <relation>_9 of the <relation>_11 of the <relation>_13
          of the <relation>_15 of the person whose <attribute_name>_17
689
         is <attribute_value>_17?
690
```

```
4. Who is the <relation>_3 of the <relation>_5 of the <relation>_7
    of the <relation>_9 of the <relation>_11 of the <relation>_13
    of the <relation> 15 of <name> 16?
```

```
693
      5. Who is the <relation>_3 of the <relation>_5 of the <relation>_7
694
          of the <relation>_9 of the <relation>_11 of the <relation>_13
695
          of the person whose <attribute_name>_15 is <attribute_value>
696
         _15?
```

```
697
      6. Who is the <relation>_3 of the <relation>_5 of the <relation>_7
698
          of the <relation>_9 of the <relation>_11 of the <relation>_13
699
          of <name>_14?
```

```
7. Who is the <relation>_3 of the <relation>_5 of the <relation>_7
700
          of the <relation>_9 of the <relation>_11 of the person whose
         <attribute_name>_13 is <attribute_value>_13?
```

702	Q	Who is the crelations 3 of the crelations 5 of the crelations 7
703	0.	of the crelations 0 of the crelations 11 of chemos 122
704	0	The the relation 2 of the relation 2 of the relation 7
705	9.	who is the cretation>_s of the cretation>_s of the cretation>_/
705		of the <relation>_9 of the person whose <attribute_name>_11</attribute_name></relation>
706		is <attribute_value>_11?</attribute_value>
707	10.	Who is the <relation>_3 of the <relation>_5 of the <relation></relation></relation></relation>
708		_/ of the <relation>_9 of <name>_10?</name></relation>
709	11.	Who is the <relation>_3 of the <relation>_5 of the <relation></relation></relation></relation>
710		_7 of the person whose <attribute_name>_9 is <attribute_value></attribute_value></attribute_name>
711		_9?
712	12.	Who is the <relation>_3 of the <relation>_5 of the <relation></relation></relation></relation>
713		_7 of <name>_8?</name>
714	13.	Who is the <relation>_3 of the <relation>_5 of the person</relation></relation>
715		whose <attribute_name>_7 is <attribute_value>_7?</attribute_value></attribute_name>
715	14.	Who is the <relation>_3 of the <relation>_5 of <name>_6?</name></relation></relation>
/16	15.	Who is the <relation>_3 of the person whose <attribute_name>_5</attribute_name></relation>
717		is <attribute_value>_5?</attribute_value>
718	16.	Who is the <relation>_3 of <name>_4?</name></relation>
719	17.	Who is the person whose <attribute_name>_3 is <attribute_value< td=""></attribute_value<></attribute_name>
720		>_3?
721	18.	What is the <attribute_name>_3 of the <relation>_4 of the <</relation></attribute_name>
722		relation>_6 of the <relation>_8 of the <relation>_10 of the <</relation></relation>
700		relation>_12 of the <relation>_14 of the <relation>_16 of the</relation></relation>
723		<relation>_18 of <name>_19?</name></relation>
724	19.	What is the <attribute name=""> 3 of the <relation> 4 of the <</relation></attribute>
725		relation> 6 of the <relation> 8 of the <relation> 10 of the <</relation></relation>
726		relation> 12 of the <relation> 14 of the <relation> 16 of the</relation></relation>
727		person whose <attribute name=""> 18 is <attribute value=""> 18?</attribute></attribute>
728	20.	What is the <attribute name=""> 3 of the <relation> 4 of the <</relation></attribute>
729	20.	relations 6 of the <relations 10="" 8="" <<="" <relations="" of="" td="" the=""></relations>
720		relation> 12 of the <relation> 14 of the <relation> 16 of <</relation></relation>
730		names 172
731	21	What is the cattribute name> 3 of the crelation> 4 of the c
732	<u> </u>	relations 6 of the crelations 8 of the crelations 10 of the c
733		relations 12 of the crelations 14 of the person whose c
734		attribute name: 16 is cattribute values 162
735	22	autipute_induce_ion is <attribute_value_io:< td=""></attribute_value_io:<>
736	22.	what is the satisfier of the selection 4 of the s
737		relation/_6 of the (relation/_6 of the (relation/_10 of the (
738	<u></u>	The inclusion of the relation of the management of the fight of the fi
700	23.	What is the <attribute_name>_3 of the <relation>_4 of the <</relation></attribute_name>
739		relation>_6 of the <relation>_8 of the <relation>_10 of the <</relation></relation>
740		relation>_12 of the person whose <attribute_name>_14 is <</attribute_name>
741	~ 4	attribute_value>_14?
742	24.	What is the <attribute_name>_3 of the <relation>_4 of the <</relation></attribute_name>
743		relation>_6 of the <relation>_8 of the <relation>_10 of the <</relation></relation>
744		relation>_12 of <name>_13?</name>
745	25.	What is the <attribute_name>_3 of the <relation>_4 of the <</relation></attribute_name>
7/6		relation>_6 of the <relation>_8 of the <relation>_10 of the</relation></relation>
747		person whose <attribute_name>_12 is <attribute_value>_12?</attribute_value></attribute_name>
747	26.	What is the <attribute_name>_3 of the <relation>_4 of the <</relation></attribute_name>
748		<pre>relation>_6 of the <relation>_8 of the <relation>_10 of <name></name></relation></relation></pre>
749		_11?
750	27.	What is the <attribute_name>_3 of the <relation>_4 of the <</relation></attribute_name>
751		relation>_6 of the <relation>_8 of the person whose <</relation>
752		attribute_name>_10 is <attribute_value>_10?</attribute_value>
753	28.	What is the <attribute_name>_3 of the <relation>_4 of the <</relation></attribute_name>
754		relation>_6 of the <relation>_8 of <name>_9?</name></relation>
755		
100		

756 29. What is the <attribute_name>_3 of the <relation>_4 of the < 757 relation>_6 of the person whose <attribute_name>_8 is <</pre> 758 attribute_value>_8? 759 30. What is the <attribute_name>_3 of the <relation>_4 of the < 760 relation>_6 of <name>_7? 761 31. What is the <attribute_name>_3 of the <relation>_4 of the person whose <attribute_name>_6 is <attribute_value>_6? 762 32. What is the <attribute_name>_3 of the <relation>_4 of <name>_5 763 ? 764 33. What is the <attribute_name>_3 of the person whose < 765 attribute_name>_4 is <attribute_value>_4? 766 34. How many <relation_plural>_2 does the <relation>_4 of the < 767 relation>_6 of the <relation>_8 of the <relation>_10 of the < 768 relation>_12 of the <relation>_14 of the <relation>_16 of the 769 <relation>_18 of <name>_19 have? 770 35. How many <relation_plural>_2 does the <relation>_4 of the < 771 relation>_6 of the <relation>_8 of the <relation>_10 of the < 772 relation>_12 of the <relation>_14 of the <relation>_16 of the person whose <attribute_name>_18 is <attribute_value>_18 have? 773 36. How many <relation_plural>_2 does the <relation>_4 of the < 774 relation>_6 of the <relation>_8 of the <relation>_10 of the < 775 relation>_12 of the <relation>_14 of the <relation>_16 of < 776 name> 17 have? 777 37. How many <relation_plural>_2 does the <relation>_4 of the < 778 relation>_6 of the <relation>_8 of the <relation>_10 of the < 779 relation>_12 of the <relation>_14 of the person whose < 780 attribute_name>_16 is <attribute_value>_16 have? 781 38. How many <relation_plural>_2 does the <relation>_4 of the < 782 relation>_6 of the <relation>_8 of the <relation>_10 of the < relation>_12 of the <relation>_14 of <name>_15 have? 783 39. How many <relation_plural>_2 does the <relation>_4 of the < 784 relation>_6 of the <relation>_8 of the <relation>_10 of the < 785 relation>_12 of the person whose <attribute_name>_14 is < 786 attribute_value>_14 have? 787 40. How many <relation_plural>_2 does the <relation>_4 of the < 788 relation>_6 of the <relation>_8 of the <relation>_10 of the < 789 relation>_12 of <name>_13 have? 790 41. How many <relation_plural>_2 does the <relation>_4 of the < 791 relation>_6 of the <relation>_8 of the <relation>_10 of the 792 person whose <attribute_name>_12 is <attribute_value>_12 have? 793 42. How many <relation_plural>_2 does the <relation>_4 of the < 794 relation>_6 of the <relation>_8 of the <relation>_10 of <name> _11 have? 795 43. How many <relation_plural>_2 does the <relation>_4 of the < 796 relation>_6 of the <relation>_8 of the person whose < 797 attribute_name>_10 is <attribute_value>_10 have? 798 44. How many <relation_plural>_2 does the <relation>_4 of the < 799 relation>_6 of the <relation>_8 of <name>_9 have? 800 45. How many <relation_plural>_2 does the <relation>_4 of the < 801 relation>_6 of the person whose <attribute_name>_8 is < 802 attribute_value>_8 have? 803 46. How many <relation_plural>_2 does the <relation>_4 of the < 804 relation>_6 of <name>_7 have? 47. How many <relation_plural>_2 does the <relation>_4 of the 805 person whose <attribute_name>_6 is <attribute_value>_6 have? 806 48. How many <relation_plural>_2 does the <relation>_4 of <name>_5 807 have? 808 49. How many <relation_plural>_2 does the person whose < 809 attribute_name>_4 is <attribute_value>_4 have?

50. How many <relation_plural>_2 does <name>_3 have?

B.3 QUESTION-ANSWER CHARACTERISTICS



Figure 5: Histogram of question difficulties (measured by reasoning steps) for universe size n = 50 at two CFG recursion limits $d \in \{10, 20\}$. We average the frequencies across 3 dataset generation seeds.





Figure 6: Distribution of number of answers across sizes $n \in \{50, 500, 5000\}$, seeds $\{1, 2, 3\}$, and CFG depth 20.

C BASELINE DETAILS

C.1 LLM SAMPLING HYPERPARAMETERS

	Temperature	Top-k	Тор-р	Repetition Penalty	Sampling Seed	Max number of output Tokens
Values	0	50	0.7	1.0	0	4096

Table 3: Default Hyperparameters values for LLM Sampling

We used the above default hyperparameters values for all models, but DeepSeek-R1-32B, where we used temperature = 0.6 and top-p = 0.95.

C.2 ZEROSHOT-SIMPLE

We use the following prompt for all models, where evidence is the concatenation of all documents in the *PhantomWiki* instance.

864 You are given the following evidence: 865 (BEGIN EVIDENCE) 866 {{evidence}} 867 (END EVIDENCE) 868 You will be provided a question. Your task is to provide an answer 869 according to these instructions: 870 - The output must be one of the following: a name (if there is 871 only one correct answer); or a list of names separated by '{ 872 constants.answer_sep}' (if there are multiple correct answers) 873 874 - DO NOT include any additional information in your answer. 875 876 Question: {{question}} 877 Answer: 878 879 For DeepSeek-R1-32B, we additionally parse the output to separate the model's reasoning process 880 from its final answer using the </think> tag. 881 882 C.3 ZEROSHOT-RAG 883 884 The prompt is exactly the same as ZEROSHOT, except we replace evidence with 4 documents 885 retrieved using the UAE-LARGE-V1. We pre-compute an index the document corpus using an FAISS 886 vector store of UAE-LARGE-V1 embeddings. Upon generation, we search for similar documents for 887 question according to maximum inner product search on document and question embeddings. 888 889 C.4 CHAIN-OF-THOUGHT-SIMPLE 890 891 We use the following prompt for all models, where evidence is replaced with a list of all documents. 892 We use a regular expression to parse the output. 893 You are given the following evidence: 894 (BEGIN EVIDENCE) 895 {{evidence}} 896 (END EVIDENCE) 897 898 You will be provided a question. Your response must end in the 899 following sentence: The answer is <answer>. 900 Here, <answer> must be one of the following: 901 - a name (if there is only one correct answer); or 902 - a list of names separated by '{constants.answer_sep}' (if there 903 are multiple correct answers). 904 Here are some examples: 905 (START OF EXAMPLES) 906 Example 1: 907 Question: Who is the brother of Dino Beltran? 908 Answer: Based on the evidence, the brother of Dino Beltran is 909 Orlando Beltran. The answer is Orlando Beltran. 910 911 Example 2: 912 Question: Who is the sibling of Barabara Beltran? 913 Answer: Based on the evidence, the siblings of Barabara Beltran 914 are Aida Wanq, Vicki Hackworth. The answer is Aida Wanq{ 915 constants.answer_sep}Vicki Hackworth. 916 Example 3: 917 Question: Who is the child of the sibling of Stacia Toombs?

918 Answer: First I need to find the sibling of Stacia Toombs. Based 919 on the evidence, the sibling of Stacia Toombs is Shelli 920 Beltran. Now I need to find the child of Shelli Beltran. Based 921 on the evidence, the children of Shelli Beltran are Aida Wang 922 , Barabara Beltran, Vicki Hackworth. The answer is Aida Wang{ 923 constants.answer_sep}Barabara Beltran{constants.answer_sep} Vicki Hackworth. 924 925 Example 4: 926 Question: Who is the uncle of William Smock? 927 Answer: An uncle is the brother of a parent. Based on the evidence 928 , the parents of William Smock are Dominique Smock, Gene Smock 929 . To find the uncle of William Smock, I need to find the 930 brother of Dominique Smock and Gene Smock. Based on the 931 evidence, Dominique Smock has no brother, and the brother of 932 Gene Smock is Eli Smock. So the uncle of William Smock is Eli 933 Smock. The answer is Eli Smock. 934 Example 5: 935 Question: What is the occupation of the sister of the grandmother 936 of Virgil Hackworth? 937 Answer: A grandmother is the mother of a parent. Based on the 938 evidence, the parents of Virgil Hackworth are Ricardo 939 Hackworth, Vicki Hackworth. To find the grandmother of Virgil 940 Hackworth, I need to find the mother of Ricardo Hackworth and 941 Vicki Hackworth. Based on the evidence, Ricardo Hackworth has 942 no mother, and the mother of Vicki Hackworth is Shelli Beltran 943 . Now I need to find the sister of Shelli Beltran. Based on 944 the evidence, the sister of Shelli Beltran is Stacia Toombs. Based on the evidence, the occupation of Stacia Toombs is 945 actuary. The answer is actuary. 946 947 Example 6: 948 Question: Who is the brother of the person whose occupation is 949 associate professor? 950 Answer: I need to search for people whose occupation is associate 951 professor. Based on the evidence, the person whose occupation 952 is associate professor is Dino Beltran. And the brother of 953 Dino Beltran is Orlando Beltran. The answer is Orlando Beltran 954 955 956 Example 7: Question: What is the date of birth of the person whose hobby is 957 meteorology? 958 Answer: I need to search for people whose hobby is meteorology. 959 Based on the evidence, the people whose hobby is meteorology 960 are Alison Smock, Barabara Beltran. The date of birth of 961 Alison Smock is 0929-10-28, and the date of birth of Barabara 962 Beltran is 0989-06-11. The answer is 0929-10-28{constants. 963 answer_sep}0989-06-11. 964 965 Example 8: 966 Question: Who is the cousin of the person whose occupation is 967 broadcast engineer? Answer: I need to search for people whose occupation is broadcast 968 engineer. Based on the evidence, the person whose occupation 969 is broadcast engineer is Barabara Beltran. A cousin is the 970 child of the sibling of the parent. Based on the evidence, the 971 parents of Barabara Beltran are Dino Beltran, Shelli Beltran.

972 The sibling of Dino Beltran is Orlando Beltran, and the 973 sibling of Shelli Beltran is Stacia Toombs. Based on the 974 evidence, Orlando Beltran has no child, and the child of 975 Stacia Toombs is Leslee Toombs. So the cousin of Barabara 976 Beltran is Leslee Toombs. The answer is Leslee Toombs. 977 Example 9: 978 Question: Who is the great-granddaughter of the person whose hobby 979 is biology? 980 Answer: I need to search for people whose hobby is biology. Based 981 on the evidence, the person whose hobby is biology is Alvaro 982 Smock. To find the great-granddaughter of Alvaro Smock, I need 983 to find the daughter of the child of the child of Alvaro 984 Smock. Based on the evidence, the children of Alvaro Smock are 985 Eli Smock, Gene Smock. Eli Smock has no child, and the child 986 of Gene Smock is Williams Smock. The daughters of Williams 987 Smock are Shelli Beltran, Stacia Toombs. So the greatgranddaughters of Alvaro Smock, whose hobby is biology, are 988 Shelli Beltran, Stacia Toombs. The answer is Shelli Beltran{ 989 constants.answer_sep}Stacia Toombs. 990 (END OF EXAMPLES) 991 992 Question: {{question}} 993 Answer: 994 995 C.5 CHAIN-OF-THOUGHT-RAG 996 997 The prompt is exactly the same as COT, except we replace evidence with 4 documents retrieved 998 using the UAE-LARGE-V1. We use the same retriever setup as described in Appendix C.3. 999 1000 C.6 REACT 1001 We modify the ReAct agent implementation of Shinn et al. (2024), which can be found at 1003 https://github.com/noahshinn/reflexion, to support our keyword-based search tools 1004 RetrieveArticle and Search. As seen in the following examples, the RetrieveArticle 1005 tool fetches the documents by title, and the Search tool fetches all document titles whose bodies 1006 contain the argument. We instruct the LLM to output predict the answer using a third tool Finish. 1007 Solve a question answering task with interleaving Thought, Action, 1008 Observation steps. 1009 Thought can reason about the current situation, and Action can be 1010 3 types: 1011 (1) RetrieveArticle[{{{entity}}}]. This action retrieves the 1012 article about {{{entity}}} if it exists. If the article does not exist, the action will say so. 1013 (2) Search[{{{attribute}}}]. This action searches the database 1014 for {{{attribute}}} and retrieves all articles that contain 1015 {{{attribute}}}. If no article contains {{{attribute}}}, 1016 the action will say so. 1017 (3) Finish[{{{answer}}}]. This action answers the question with 1018 $\{\{\{answer\}\}\}\}.$ 1019 If you cannot find the answer, output the empty answer like: 1020 Finish[]. 1021 If there are multiple answers A,B,C, answer with a list like: 1022 Finish[A{constants.answer_sep}B{constants.answer_sep}C]. 1023 You may take as many steps as necessary. 1024 Here are some examples: 1025 (START OF EXAMPLES)

1026 Example 1: 1027 Question: Who is the brother of Dino Beltran? 1028 Thought 1: I need to retrieve article about Dino Beltran and find 1029 who his brother is. 1030 Action 1: RetrieveArticle[Dino Beltran]. Observation 1: # Dino Beltran ## Family The parents of Dino 1031 Beltran are Brian Beltran, Daisy Beltran. Dino Beltran's 1032 sibling is Orlando Beltran. The brother of Dino Beltran is 1033 Orlando Beltran. The mother of Dino Beltran is Daisy Beltran. 1034 The father of Dino Beltran is Brian Beltran. The children 1035 of Dino Beltran are Aida Wang, Barabara Beltran, Vicki 1036 Hackworth. The daughters of Dino Beltran are Aida Wang, 1037 Barabara Beltran, Vicki Hackworth. The wife of Dino Beltran 1038 is Shelli Beltran. ## Friends The friend of Dino Beltran is 1039 Alvaro Smock. ## Attributes The date of birth of Dino Beltran 1040 is 0958-08-09. The occupation of Dino Beltran is associate professor. The hobby of Dino Beltran is shoqi. 1041 1042 Thought 2: The brother of Dino Beltran is Orlando Beltran, so the answer is Orlando Beltran. 1043 Action 2: Finish[Orlando Beltran]. 1044 1045 Example 2: 1046 Question: Who is the mother of Ivana Smith? 1047 Thought 1: I need to retrieve article about Ivana Smith and find 1048 who her mother is. 1049 Action 1: RetrieveArticle[Ivana Smith]. 1050 Observation 1: No article exists for the requested entity. Please 1051 try retrieving article for another entity. 1052 Thought 2: Article about Ivana Smith does not exist, so I cannot find the mother of Ivana Smith. I will output empty answer. 1053 Action 2: Finish[]. 1054 1055 Example 3: 1056 Question: Who is the sibling of Barabara Beltran? 1057 Thought 1: I need to retrieve article about Barabara Beltran and 1058 find who her sibling is. 1059 Action 1: RetrieveArticle[Barabara Beltran]. 1060 Observation 1: # Barabara Beltran ## Family The parents of 1061 Barabara Beltran are Dino Beltran, Shelli Beltran. Barabara 1062 Beltran's siblings are Aida Wang, Vicki Hackworth. The 1063 sisters of Barabara Beltran are Aida Wang, Vicki Hackworth. 1064 The mother of Barabara Beltran is Shelli Beltran. The father of Barabara Beltran is Dino Beltran. ## Friends ## Attributes 1065 The date of birth of Barabara Beltran is 0989-06-11. The 1066 occupation of Barabara Beltran is broadcast engineer. The 1067 hobby of Barabara Beltran is meteorology. 1068 Thought 2: The siblings of Barabara Beltran are Aida Wang, Vicki 1069 Hackworth. So the answer is Aida Wang, Vicki Hackworth. 1070 Action 2: Finish[Aida Wang{constants.answer_sep}Vicki Hackworth]. 1071 1072 Example 4: 1073 Question: Who is the child of the sibling of Stacia Toombs? 1074 Thought 1: First, I need to retrieve article about Stacia Toombs 1075 and find who her sibling is. Action 1: RetrieveArticle[Stacia Toombs]. 1076 Observation 1: # Stacia Toombs ## Family The parents of Stacia 1077 Toombs are Alison Smock, Williams Smock. Stacia Toombs's 1078 sibling is Shelli Beltran. The sister of Stacia Toombs is 1079 Shelli Beltran. The mother of Stacia Toombs is Alison Smock.

1080 The father of Stacia Toombs is Williams Smock. The child of 1081 Stacia Toombs is Leslee Toombs. The daughter of Stacia Toombs 1082 is Leslee Toombs. The husband of Stacia Toombs is Wilbert 1083 Toombs. ## Friends The friends of Stacia Toombs are Brian 1084 Beltran, Isiah Lutz, Leeann Hackworth, Lesley Lutz, Ryan Wang. 1085 ## Attributes The date of birth of Stacia Toombs is 0959-03-22. The occupation of Stacia Toombs is actuary. The 1086 hobby of Stacia Toombs is finance. 1087 Thought 2: The sibling of Stacia Toombs is Shelli Beltran. Now I 1088 need to retrieve article about Shelli Beltran and find who her 1089 child is. 1090 Action 2: RetrieveArticle[Shelli Beltran]. 1091 Observation 2: # Shelli Beltran ## Family The parents of Shelli 1092 Beltran are Alison Smock, Williams Smock. Shelli Beltran's 1093 sibling is Stacia Toombs. The sister of Shelli Beltran is 1094 Stacia Toombs. The mother of Shelli Beltran is Alison Smock. 1095 The father of Shelli Beltran is Williams Smock. The children 1096 of Shelli Beltran are Aida Wang, Barabara Beltran, Vicki Hackworth. The daughters of Shelli Beltran are Aida Wang, 1097 Barabara Beltran, Vicki Hackworth. The husband of Shelli 1098 Beltran is Dino Beltran. ## Friends The friends of Shelli 1099 Beltran are Brian Beltran, Eli Smock, Isiah Lutz, Leslee 1100 Toombs, Lesley Lutz, Ryan Wang. ## Attributes The date of 1101 birth of Shelli Beltran is 0958-03-08. The occupation of 1102 Shelli Beltran is occupational therapist. The hobby of Shelli 1103 Beltran is sociology. 1104 Thought 3: The child of Shelli Beltran is Aida Wang, Barabara 1105 Beltran, Vicki Hackworth. So the answer is Aida Wang, Barabara 1106 Beltran, Vicki Hackworth. Action 3: Finish[Aida Wang{constants.answer_sep}Barabara Beltran{ 1107 constants.answer_sep}Vicki Hackworth]. 1108 1109 Example 5: 1110 Question: Who is the uncle of William Smock? 1111 Thought 1: First, I need to retrieve article about William Smock. 1112 To find the uncle of William Smock, I need to find the brother 1113 of his parent. 1114 Action 1: RetrieveArticle[William Smock]. 1115 Observation 1: # Williams Smock ## Family The parents of Williams 1116 Smock are Dominique Smock, Gene Smock. The mother of Williams 1117 Smock is Dominique Smock. The father of Williams Smock is Gene Smock. The children of Williams Smock are Shelli Beltran 1118 , Stacia Toombs. The daughters of Williams Smock are Shelli 1119 Beltran, Stacia Toombs. The wife of Williams Smock is Alison 1120 Smock. ## Friends The friend of Williams Smock is Lannie 1121 ## Attributes The date of birth of Williams Smock is Smock. 1122 0926-04-04. The occupation of Williams Smock is clinical 1123 biochemist. The hobby of Williams Smock is social studies. 1124 Thought 2: The parents of William Smock are Dominique Smock, Gene 1125 Smock. To find the uncle of William Smock, I need to find the 1126 brother of Dominique Smock and Gene Smock. 1127 Action 2: RetrieveArticle[Dominique Smock]. 1128 Observation 2: # Dominique Smock ## Family The child of Dominique Smock is Williams Smock. The son of Dominique Smock is 1129 Williams Smock. The husband of Dominique Smock is Gene Smock. 1130 ## Friends The friends of Dominique Smock are Gene Smock, 1131 Isiah Lutz, Orlando Beltran, Vicki Hackworth. ## Attributes 1132 The date of birth of Dominique Smock is 0897-09-08. The 1133

```
1134
         occupation of Dominique Smock is sports therapist. The hobby
1135
         of Dominique Smock is dominoes.
1136
      Thought 3: Dominique Smock has no brother. Now I need to retrieve
1137
         article about Gene Smock and find his brother.
1138
      Action 3: RetrieveArticle[Gene Smock].
      Observation 3: # Gene Smock ## Family The parents of Gene Smock
1139
         are Alvaro Smock, Lannie Smock. Gene Smock's sibling is Eli
1140
         Smock. The brother of Gene Smock is Eli Smock. The mother of
1141
          Gene Smock is Lannie Smock. The father of Gene Smock is
1142
         Alvaro Smock. The child of Gene Smock is Williams Smock.
                                                                      The
1143
          son of Gene Smock is Williams Smock. The wife of Gene Smock
1144
         is Dominique Smock. ## Friends The friends of Gene Smock are
1145
         Leeann Hackworth, Leisa Lutz, Ricardo Hackworth, Alvaro Smock,
1146
          Dominique Smock. ## Attributes The date of birth of Gene
1147
         Smock is 0898-08-16. The occupation of Gene Smock is
1148
         immunologist. The hobby of Gene Smock is architecture.
1149
      Thought 4: The brother of Gene Smock is Eli Smock, and Dominique
1150
         Smock has no brother. Since Gene Smock and Dominique Smock are
          the parents of William Smock, the uncle of William Smock is
1151
         Eli Smock. So the answer is Eli Smock.
1152
      Action 4: Finish[Eli Smock].
1153
1154
      Example 6:
1155
      Question: What is the occupation of the sister of the grandmother
1156
         of Virgil Hackworth?
1157
      Thought 1: First, I need to retrieve article about Virgil
1158
         Hackworth. To find the grandmother of Virgil Hackworth, I need
1159
          to find the mother of his parent.
1160
      Action 1: RetrieveArticle[Virgil Hackworth].
      Observation 1: # Virgil Hackworth ## Family The parents of Virgil
1161
         Hackworth are Ricardo Hackworth, Vicki Hackworth. Virgil
1162
         Hackworth's siblings are Leeann Hackworth, Leisa Lutz. The
1163
         sisters of Virgil Hackworth are Leeann Hackworth, Leisa Lutz.
1164
          The mother of Virgil Hackworth is Vicki Hackworth. The
1165
         father of Virgil Hackworth is Ricardo Hackworth. ## Friends
1166
         The friends of Virgil Hackworth are Alison Smock, Eli Smock,
1167
         Orlando Beltran, Ryan Wang. ## Attributes The date of birth
1168
         of Virgil Hackworth is 1009-06-23. The occupation of Virgil
1169
         Hackworth is sound technician. The hobby of Virgil Hackworth
1170
         is wikipedia editing.
1171
      Thought 2: The parents of Virgil Hackworth are Ricardo Hackworth,
1172
         Vicki Hackworth. To find the grandmother of Virgil Hackworth,
         I need to find the mother of Ricardo Hackworth and Vicki
1173
         Hackworth.
1174
      Action 2: RetrieveArticle[Ricardo Hackworth].
1175
      Observation 2: # Ricardo Hackworth ## Family The children of
1176
         Ricardo Hackworth are Leeann Hackworth, Leisa Lutz, Virgil
1177
         Hackworth. The son of Ricardo Hackworth is Virgil Hackworth.
1178
          The daughters of Ricardo Hackworth are Leeann Hackworth,
1179
         Leisa Lutz. The wife of Ricardo Hackworth is Vicki Hackworth.
1180
           ## Friends The friends of Ricardo Hackworth are Gene Smock,
1181
         Isiah Lutz, Johnetta Wang. ## Attributes The date of birth of
1182
          Ricardo Hackworth is 0983-02-24. The occupation of Ricardo
         Hackworth is clinical research associate. The hobby of
1183
         Ricardo Hackworth is dairy farming.
1184
      Thought 3: Ricardo Hackworth has no mother. So I need to retrieve
1185
         article about Vicki Hackworth and find her mother.
1186
      Action 3: RetrieveArticle[Vicki Hackworth].
1187
```

1188 Observation 3: # Vicki Hackworth ## Family The parents of Vicki 1189 Hackworth are Dino Beltran, Shelli Beltran. Vicki Hackworth's 1190 siblings are Aida Wang, Barabara Beltran. The sisters of 1191 Vicki Hackworth are Aida Wang, Barabara Beltran. The mother 1192 of Vicki Hackworth is Shelli Beltran. The father of Vicki Hackworth is Dino Beltran. The children of Vicki Hackworth 1193 are Leeann Hackworth, Leisa Lutz, Virgil Hackworth. The son 1194 of Vicki Hackworth is Virgil Hackworth. The daughters of 1195 Vicki Hackworth are Leeann Hackworth, Leisa Lutz. The husband 1196 of Vicki Hackworth is Ricardo Hackworth. ## Friends The 1197 friends of Vicki Hackworth are Brian Beltran, Dominique Smock, 1198 Eli Smock. ## Attributes The date of birth of Vicki 1199 Hackworth is 0985-05-30. The occupation of Vicki Hackworth is 1200 police officer. The hobby of Vicki Hackworth is meditation. 1201 Thought 4: The mother of Vicki Hackworth is Shelli Beltran. So the 1202 only grandmother of Virgil Hackworth is Shelli Beltran. Now I 1203 need to retrieve article about Shelli Beltran and find her 1204 sister. Action 4: RetrieveArticle[Shelli Beltran]. 1205 Observation 4: # Shelli Beltran ## Family The parents of Shelli 1206 Beltran are Alison Smock, Williams Smock. Shelli Beltran's 1207 sibling is Stacia Toombs. The sister of Shelli Beltran is 1208 Stacia Toombs. The mother of Shelli Beltran is Alison Smock. 1209 The father of Shelli Beltran is Williams Smock. The children 1210 of Shelli Beltran are Aida Wang, Barabara Beltran, Vicki 1211 Hackworth. The daughters of Shelli Beltran are Aida Wang, 1212 Barabara Beltran, Vicki Hackworth. The husband of Shelli 1213 Beltran is Dino Beltran. ## Friends The friends of Shelli 1214 Beltran are Brian Beltran, Eli Smock, Isiah Lutz, Leslee Toombs, Lesley Lutz, Ryan Wang. ## Attributes The date of 1215 birth of Shelli Beltran is 0958-03-08. The occupation of 1216 Shelli Beltran is occupational therapist. The hobby of Shelli 1217 Beltran is sociology. 1218 Thought 5: The sister of Shelli Beltran is Stacia Toombs. Now I 1219 need to retrieve article about Stacia Toombs and find her 1220 occupation. 1221 Action 5: RetrieveArticle[Stacia Toombs]. 1222 Observation 5: # Stacia Toombs ## Family The parents of Stacia 1223 Toombs are Alison Smock, Williams Smock. Stacia Toombs's 1224 sibling is Shelli Beltran. The sister of Stacia Toombs is 1225 Shelli Beltran. The mother of Stacia Toombs is Alison Smock. 1226 The father of Stacia Toombs is Williams Smock. The child of Stacia Toombs is Leslee Toombs. The daughter of Stacia Toombs 1227 is Leslee Toombs. The husband of Stacia Toombs is Wilbert 1228 Toombs. ## Friends The friends of Stacia Toombs are Brian 1229 Beltran, Isiah Lutz, Leeann Hackworth, Lesley Lutz, Ryan Wang. 1230 ## Attributes The date of birth of Stacia Toombs is 1231 0959-03-22. The occupation of Stacia Toombs is actuary. The 1232 hobby of Stacia Toombs is finance. 1233 Thought 6: The occupation of Stacia Toombs is actuary. So the 1234 answer is actuary. 1235 Action 6: Finish[actuary]. 1236 1237 Example 7: Question: Who is the brother of the person whose occupation is 1238 associate professor? 1239 Thought 1: First, I need to search for people whose occupation is 1240 associate professor. 1241 Action 1: Search[associate professor].

1242 Observation 1: (1) Dino Beltran 1243 Thought 2: Person whose hobby is associate professor is Dino 1244 Beltran. Now I need to retrieve article about Dino Beltran and 1245 find his brother. 1246 Action 2: RetrieveArticle[Dino Beltran]. Observation 2: # Dino Beltran ## Family The parents of Dino 1247 Beltran are Brian Beltran, Daisy Beltran. Dino Beltran's 1248 sibling is Orlando Beltran. The brother of Dino Beltran is 1249 Orlando Beltran. The mother of Dino Beltran is Daisy Beltran. 1250 The father of Dino Beltran is Brian Beltran. The children 1251 of Dino Beltran are Aida Wang, Barabara Beltran, Vicki 1252 Hackworth. The daughters of Dino Beltran are Aida Wang, 1253 Barabara Beltran, Vicki Hackworth. The wife of Dino Beltran 1254 is Shelli Beltran. ## Friends The friend of Dino Beltran is 1255 Alvaro Smock. ## Attributes The date of birth of Dino Beltran 1256 is 0958-08-09. The occupation of Dino Beltran is associate professor. The hobby of Dino Beltran is shogi. 1257 1258 Thought 3: The brother of Dino Beltran is Orlando Beltran, so the answer is Orlando Beltran. 1259 Action 3: Finish[Orlando Beltran]. 1260 1261 Example 8: 1262 Question: What is the date of birth of the person whose hobby is 1263 meteorology? 1264 Thought 1: First, I need to search for people whose hobby is 1265 meteorology. 1266 Action 1: Search[meteorology]. 1267 Observation 1: (1) Alison Smock (2) Barabara Beltran Thought 2: People whose hobby is meteorology are Alison Smock, 1268 Barabara Beltran. Now I need to retrieve article about Alison 1269 Smock and Barabara Beltran and find their date of birth. 1270 Action 2: RetrieveArticle[Alison Smock]. 1271 Observation 2: # Alison Smock ## Family The children of Alison 1272 Smock are Shelli Beltran, Stacia Toombs. The daughters of 1273 Alison Smock are Shelli Beltran, Stacia Toombs. The husband 1274 of Alison Smock is Williams Smock. ## Friends The friends of 1275 Alison Smock are Eli Smock, Lannie Smock, Virgil Hackworth. 1276 ## Attributes The date of birth of Alison Smock is 0929-10-28. 1277 The occupation of Alison Smock is health promotion 1278 specialist. The hobby of Alison Smock is meteorology. 1279 Thought 3: The date of birth of Alison Smock is 0929-10-28. Now I 1280 need to retrieve article about Barabara Beltran and find her date of birth. 1281 Action 3: RetrieveArticle[Barabara Beltran]. 1282 Observation 3: # Barabara Beltran ## Family The parents of 1283 Barabara Beltran are Dino Beltran, Shelli Beltran. Barabara 1284 Beltran's siblings are Aida Wang, Vicki Hackworth. The 1285 sisters of Barabara Beltran are Aida Wang, Vicki Hackworth. 1286 The mother of Barabara Beltran is Shelli Beltran. The father 1287 of Barabara Beltran is Dino Beltran. ## Friends ## Attributes 1288 The date of birth of Barabara Beltran is 0989-06-11. The 1289 occupation of Barabara Beltran is broadcast engineer. The 1290 hobby of Barabara Beltran is meteorology. Thought 4: The date of birth of Barabara Beltran is 0989-06-11. So 1291 the answer is 0929-10-28, 0989-06-11. 1292 Action 4: Finish[0929-10-28{constants.answer sep}0989-06-11]. 1293 1294 Example 9: 1295

1296 Question: Who is the cousin of the person whose occupation is 1297 broadcast engineer? 1298 Thought 1: First, I need to search for people whose occupation is 1299 broadcast engineer. 1300 Action 1: Search[broadcast engineer]. Observation 1: # Barabara Beltran ## Family The parents of 1301 Barabara Beltran are Dino Beltran, Shelli Beltran. Barabara 1302 Beltran's siblings are Aida Wang, Vicki Hackworth. The 1303 sisters of Barabara Beltran are Aida Wang, Vicki Hackworth. 1304 The mother of Barabara Beltran is Shelli Beltran. The father 1305 of Barabara Beltran is Dino Beltran. ## Friends ## Attributes 1306 The date of birth of Barabara Beltran is 0989-06-11. The 1307 occupation of Barabara Beltran is broadcast engineer. The 1308 hobby of Barabara Beltran is meteorology. 1309 Thought 2: Person whose hobby is broadcast engineer is Barabara 1310 Beltran. The cousin of Barabara Beltran is the child of 1311 sibling of her parent. Barabara Beltran's parents are Dino 1312 Beltran and Shelly Beltran. Now I need to find the sibling of Dino Beltran and Shelly Beltran. 1313 Action 2: RetrieveArticle[Dino Beltran]. 1314 Observation 2: # Dino Beltran ## Family The parents of Dino 1315 Beltran are Brian Beltran, Daisy Beltran. Dino Beltran's 1316 sibling is Orlando Beltran. The brother of Dino Beltran is 1317 Orlando Beltran. The mother of Dino Beltran is Daisy Beltran. 1318 The father of Dino Beltran is Brian Beltran. The children 1319 of Dino Beltran are Aida Wang, Barabara Beltran, Vicki 1320 Hackworth. The daughters of Dino Beltran are Aida Wang, 1321 Barabara Beltran, Vicki Hackworth. The wife of Dino Beltran 1322 is Shelli Beltran. ## Friends The friend of Dino Beltran is Alvaro Smock. ## Attributes The date of birth of Dino Beltran 1323 is 0958-08-09. The occupation of Dino Beltran is associate 1324 professor. The hobby of Dino Beltran is shogi. 1325 Thought 3: The sibling of Dino Beltran is Orlando Beltran. Now I 1326 need to retrieve article about Orlando Beltran and find his 1327 child. 1328 Action 3: RetrieveArticle[Orlando Beltran]. 1329 Observation 3: # Orlando Beltran ## Family The parents of Orlando 1330 Beltran are Brian Beltran, Daisy Beltran. Orlando Beltran's 1331 sibling is Dino Beltran. The brother of Orlando Beltran is 1332 Dino Beltran. The mother of Orlando Beltran is Daisy Beltran. 1333 The father of Orlando Beltran is Brian Beltran. ## Friends 1334 The friends of Orlando Beltran are Virgil Hackworth, Dominique ## Attributes The date of birth of Orlando Beltran is Smock. 1335 0953-10-23. The occupation of Orlando Beltran is petroleum 1336 engineer. The hobby of Orlando Beltran is learning. 1337 Thought 4: Orlando Beltran has no child, so no cousin from the 1338 side of Dino Beltran. Now I need to retrieve article about 1339 Shelli Beltran and find her sibling. 1340 Action 4: RetrieveArticle[Shelli Beltran]. 1341 Observation 4: # Shelli Beltran ## Family The parents of Shelli 1342 Beltran are Alison Smock, Williams Smock. Shelli Beltran's 1343 sibling is Stacia Toombs. The sister of Shelli Beltran is 1344 Stacia Toombs. The mother of Shelli Beltran is Alison Smock. The father of Shelli Beltran is Williams Smock. The children 1345 of Shelli Beltran are Aida Wang, Barabara Beltran, Vicki 1346 Hackworth. The daughters of Shelli Beltran are Aida Wang, 1347 Barabara Beltran, Vicki Hackworth. The husband of Shelli 1348 Beltran is Dino Beltran. ## Friends The friends of Shelli 1349 Beltran are Brian Beltran, Eli Smock, Isiah Lutz, Leslee

1350 Toombs, Lesley Lutz, Ryan Wang. ## Attributes The date of 1351 birth of Shelli Beltran is 0958-03-08. The occupation of 1352 Shelli Beltran is occupational therapist. The hobby of Shelli 1353 Beltran is sociology. 1354 Thought 5: The sibling of Shelli Beltran is Stacia Toombs. Now I need to retrieve article about Stacia Toombs and find her 1355 child, to get the cousin of Barabara Beltran. 1356 Action 5: RetrieveArticle[Stacia Toombs]. 1357 Observation 5: # Stacia Toombs ## Family The parents of Stacia 1358 Toombs are Alison Smock, Williams Smock. Stacia Toombs's 1359 sibling is Shelli Beltran. The sister of Stacia Toombs is 1360 Shelli Beltran. The mother of Stacia Toombs is Alison Smock. 1361 The father of Stacia Toombs is Williams Smock. The child of 1362 The daughter of Stacia Toombs Stacia Toombs is Leslee Toombs. 1363 is Leslee Toombs. The husband of Stacia Toombs is Wilbert 1364 Toombs. ## Friends The friends of Stacia Toombs are Brian 1365 Beltran, Isiah Lutz, Leeann Hackworth, Lesley Lutz, Ryan Wang. ## Attributes The date of birth of Stacia Toombs is 1366 0959-03-22. The occupation of Stacia Toombs is actuary. 1367 The hobby of Stacia Toombs is finance. 1368 Thought 6: The child of Stacia Toombs is Leslee Toombs. So the 1369 cousin of Barabara Beltran is Leslee Toombs. The answer is 1370 Leslee Toombs. 1371 Action 6: Finish[Leslee Toombs]. 1372 1373 Example 10: 1374 Question: Who is the great-granddaughter of the person whose hobby 1375 is biology? 1376 Thought 1: First, I need to search for people whose hobby is 1377 biology. Action 1: Search[biology]. 1378 Observation 1: (1) Alvaro Smock 1379 Thought 2: Person whose hobby is biology is Alvaro Smock. Now I 1380 need to retrieve article about Alvaro Smock and find his great 1381 -granddaughter. A great-granddaughter is the daughter of the 1382 child of the child of the person. 1383 Action 2: RetrieveArticle[Alvaro Smock]. 1384 Observation 2: # Alvaro Smock ## Family The children of Alvaro 1385 Smock are Eli Smock, Gene Smock. The sons of Alvaro Smock are 1386 Eli Smock, Gene Smock. The wife of Alvaro Smock is Lannie 1387 Smock. ## Friends The friends of Alvaro Smock are Dino 1388 Beltran, Gene Smock, Aida Wang. ## Attributes The date of birth of Alvaro Smock is 0867-07-12. The occupation of Alvaro 1389 Smock is osteopath. The hobby of Alvaro Smock is biology. 1390 Thought 3: The children of Alvaro Smock are Eli Smock, Gene Smock. 1391 Now I need to retrieve article about Eli Smock and Gene Smock 1392 and find their child. 1393 Action 3: RetrieveArticle[Eli Smock]. 1394 Observation 3: # Eli Smock ## Family The parents of Eli Smock are 1395 Alvaro Smock, Lannie Smock. Eli Smock's sibling is Gene Smock 1396 The brother of Eli Smock is Gene Smock. The mother of Eli 1397 Smock is Lannie Smock. The father of Eli Smock is Alvaro 1398 ## Friends The friends of Eli Smock are Leisa Lutz, Smock. Shelli Beltran, Vicki Hackworth, Virgil Hackworth, Alison 1399 Smock, Brian Beltran. ## Attributes The date of birth of Eli 1400 Smock is 0901-01-18. The occupation of Eli Smock is retail 1401 manager. The hobby of Eli Smock is tether car. 1402 Thought 4: Eli Smock has no child. Now I need to retrieve article 1403 about Gene Smock and find his child.

